

(19) World Intellectual Property Organization  
International Bureau



(43) International Publication Date  
25 April 2002 (25.04.2002)

PCT

(10) International Publication Number  
**WO 02/33915 A1**

(51) International Patent Classification<sup>7</sup>: **H04L 12/56**

(21) International Application Number: PCT/US01/31259

(22) International Filing Date: 5 October 2001 (05.10.2001)

(25) Filing Language: English

(26) Publication Language: English

(30) Priority Data:

|            |                                |    |
|------------|--------------------------------|----|
| 60/241,450 | 17 October 2000 (17.10.2000)   | US |
| 60/275,206 | 12 March 2001 (12.03.2001)     | US |
| 09/903,441 | 10 July 2001 (10.07.2001)      | US |
| 09/903,423 | 10 July 2001 (10.07.2001)      | US |
| 09/923,924 | 6 August 2001 (06.08.2001)     | US |
| 09/960,623 | 20 September 2001 (20.09.2001) | US |

(63) Related by continuation (CON) or continuation-in-part (CIP) to earlier applications:

|          |                                |
|----------|--------------------------------|
| US       | 09/960,623 (CIP)               |
| Filed on | 20 September 2001 (20.09.2001) |
| US       | 09/923,924 (CIP)               |
| Filed on | 6 August 2001 (06.08.2001)     |

|          |                              |
|----------|------------------------------|
| US       | 09/903,441 (CIP)             |
| Filed on | 10 July 2001 (10.07.2001)    |
| US       | 09/903,423 (CIP)             |
| Filed on | 10 July 2001 (10.07.2001)    |
| US       | 60/275,206 (CIP)             |
| Filed on | 12 March 2001 (12.03.2001)   |
| US       | 60/241,450 (CIP)             |
| Filed on | 17 October 2000 (17.10.2000) |

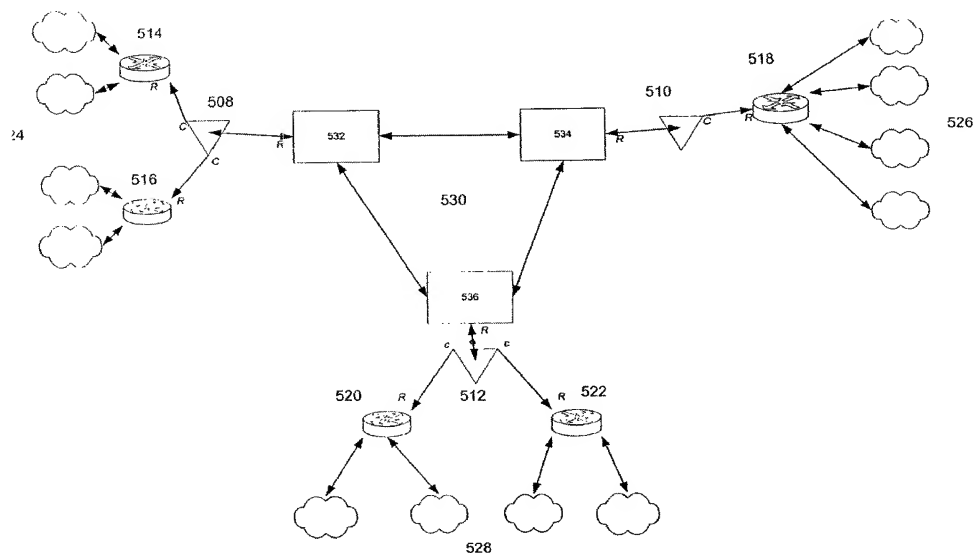
(71) Applicant (for all designated States except US): **ROUTE-SCIENCE TECHNOLOGIES INC** [US/US]; 167 2nd Avenue, San Mateo, CA 94401 (US).

(72) Inventors; and

(75) Inventors/Applicants (for US only): **BALDONADO, Omar, C** [US/US]; 700 Alester Avenue, Palo Alto, CA 94303 (US). **FINN, Sean, P** [US/US]; 1533 Escondido Way, Belmont, CA 94002 (US). **KARAM, Mansour, J.** [LB/US]; #421, 707 Continental Circle, Mountain View, CA 94040 (US). **LLOYD, Michael, A** [US/US]; 160 Arundel Road, San Carlos, CA 94070 (US). **MADAN, Herbert, S.** [IN/US]; 347 Blackfield Drive, Tiburon, CA 94920 (US). **McGUIRE, James, G** [US/US]; 2312 Gough

[Continued on next page]

(54) Title: METHOD AND APPARATUS FOR COORDINATING ROUTING PARAMETERS VIA A BACK-CHANNEL COMMUNICATION MEDIUM



(57) Abstract: Systems and methods are described for enabling routers to coordinate via a back-channel communication medium. The information exchanged over the back-channel is used to increase the number of paths considered for the routers during route optimization. The Decision Makers may assert routes and prefixes to the routers under their control. This may be done via a Border Gateway Protocol (BGP) feed. The Decision Makers, in turn, communicate separately with one another, in order to coordinate routing policy amongst themselves. This coordination may be performed over a back-channel, which may take the form of physical or logical connections between the Decision Makers.



WO 02/33915 A1



Street, San Francisco, CA 94019 (US). **VILLAYERDE, Jose-Miguel, Pulido** [ES/US]; 1020 Bryant Street, Palo Alto, CA 94301 (US).

(74) **Agent:** **SUZUE, Kenta**; Wilson Sonsini Goodrich & Rosati, 650 Page Mill Road, Palo Alto, CA 94304-1050 (US).

(81) **Designated States (national):** AE, AG, AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, BZ, CA, CH, CN, CO, CR, CU, CZ, DE, DK, DM, DZ, EC, EE, ES, FI, GB, GD, GE, GH, GM, HR, HU, ID, IL, IN, IS, JP, KE, KG, KP, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MA, MD, MG, MK, MN, MW, MX, MZ, NO, NZ, PH, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, TZ, UA, UG, US, UZ, VN, YU, ZA, ZW.

(84) **Designated States (regional):** ARIPO patent (GH, GM, KE, LS, MW, MZ, SD, SL, SZ, TZ, UG, ZW), Eurasian patent (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), European patent (AT, BE, CH, CY, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE, TR), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, ML, MR, NE, SN, TD, TG).

**Published:**

- with international search report
- before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments

*For two-letter codes and other abbreviations, refer to the "Guidance Notes on Codes and Abbreviations" appearing at the beginning of each regular issue of the PCT Gazette.*

**METHOD AND APPARATUS FOR COORDINATING ROUTING  
PARAMETERS VIA A BACK-CHANNEL COMMUNICATION  
MEDIUM**

**BACKGROUND OF THE INVENTION**

**5      Field of the Invention**

This invention relates to the field of networking. In particular, the invention relates to systems and methods for coordinating routing information amongst routers.

**Description of the Related Art**

10            Internetworks such as the Internet are currently comprised of Autonomous Systems, which exchange routing information via exterior gateway protocols. Amongst the most important of these protocols is the Border Gateway Protocol, or BGP. BGPv4 constructs a directed graph of the Autonomous Systems, based on the information exchanged between BGP  
15 routers. Each Autonomous System is identified by a unique 16 bit AS number, and BGP ensures loop-free routing amongst the Autonomous Systems; BGP also enables the exchange of additional routing information between Autonomous Systems. BGP is further described in several RFCs, which are compiled in The Big Book of Border Gateway Protocol RFCs, by Pete Loshin,  
20 which is hereby incorporated by reference.

The Border Gateway Protocol provides network administrators some measure of control over outbound traffic control from their respective organizations. For instance, the protocol includes a LOCAL\_PREF attribute, which allows BGP speakers to inform other BGP speakers within the  
25 Autonomous System of the speaker's preference for an advertised route. The local preference attribute includes a degree of preference for the advertised route, which enables comparison against other routes for the same destination. As the LOCAL\_PREF attribute is shared with other routers within an Autonomous System via IBGP, it determines outbound routes used by routers  
30 within the Autonomous System.

A WEIGHT parameter may also be used to indicate route preferences; higher preferences are assigned to routes with higher values of WEIGHT. The

WEIGHT parameter is a proprietary addition to the BGPv4 supported by Cisco Systems, Inc. of San Jose, CA. In typical implementations, the WEIGHT parameter is given higher precedence than other BGP attributes.

5 The performance knobs described above are, however, rather simple, as they do not offer system administrators with sufficiently sophisticated means for enabling routers to discriminate amongst routes. There is a need for technology that enables greater control over outbound routing policy. In particular, there is a need to allow performance data about routes to be exchanged between routers. Additionally, system administrators should be able to fine tune routing policy  
10 based upon sophisticated, up-to-date measurements of route performance and pricing analysis of various routes.

### SUMMARY OF THE INVENTION

The invention includes systems and methods for enabling networking devices to coordinate via a back-channel communication medium. The  
15 information exchanged over the back-channel is used to increase the number of paths considered for the routers during route optimization.

In embodiments of the invention, a set of Routing Intelligence Units may be used to control a set of routers, such that each Routing Intelligence Unit controls a distinct subset of the routers. The Routing Intelligence Units may  
20 assert routes to the routers under their control. In some embodiments, this is done via a Border Gateway Protocol (BGP) feed. The Decision Makers, in turn, communicate separately with one another, in order to coordinate routing policy amongst themselves. This coordination may be performed over a back-channel, which may take the form of physical or logical connections between the  
25 Routing Intelligence Units. In some embodiments, communications over the back-channel are conducted via separate BGP sessions. In embodiments utilizing BGP for communication to the routers and the back-channel, the Routing Intelligence Unit may be configured as a route-reflector client to both other decision makers and the routers it controls. This ensures that the Routing  
30 Intelligence Unit does not simply transmit information in either direction without consideration.

In some embodiments of the invention, a Routing Intelligence Unit send

updates to other Routing Intelligence Units whenever the Routing Intelligence Unit is also asserting to the routers under its control. In alternative embodiments, the Routing Intelligence Unit may send updates when it decides that the current routes are correct.

5           In some embodiments of the invention, performance scores for prefixes are communicated between Routing Intelligence Units. In some of the embodiments utilizing BGP for such coordination, these performance scores are translated to units of Local Preference. This ensures that the Routing Intelligence Units will automatically select and propagate the best score.

10           Some embodiments of the invention include techniques enabling Routing Intelligence Units to evaluate prefixes that arrive via coordination. In some embodiments, when local and remote routes have comparable scores, the local route is chosen by default. In other embodiments, a static penalty is applied to all remote announcements. In some embodiments, dynamic penalties  
15           are applied. These and other embodiments are described in greater detail infra.

## BRIEF DESCRIPTION OF THE FIGURES

Fig. 1 – Fig.4 illustrate different configurations of routing intelligence units and edge routers, according to some embodiments of the invention.

20           Figure 5a schematically illustrates an internal architecture of a routing intelligence unit according to some embodiments of the invention.

Figure 5b illustrates coordination between routing intelligence units via a back-channel according to embodiments of the invention.

Figure 6 illustrates a queuing and threading structure used in the routing intelligence unit in some embodiments of the invention.

## 25           DETAILED DESCRIPTION

### A. System Overview

In some embodiments of the invention, one or more routing intelligence units are stationed at the premises of a multi-homed organization, each of which controls one or more edge routers. These devices inject BGP updates to the

Edge Routers they control, based on performance data from measurements obtained locally, or from a Routing Intelligence Exchange—Routing Intelligence Exchanges are further described in U.S. Provisional Applications No. 60/241,450, filed October 17, 2000 and U.S. Provisional Application No. 60/275,206, filed March 12, 2001, and U.S. Applications No. 09/903,441, filed July 10, 2001, U.S. Application No. 09/923,924, filed August 6, 2001, and U.S. Application No. 09/903,423, filed July 10, 2001, which are hereby incorporated by reference in their entirety. Different configurations of these routing intelligence units and edge routers are illustrated in Figures 1 through 4. In some embodiments illustrated in Figure 1, one edge router 102 with multiple ISPs 104 and 106 is controlled by a single device 100. Figure 2 illustrates embodiments in which the routing intelligence unit 200 controls multiple edge routers 202 and 204, each of which in turn links to multiple ISPs 206, 208, 210, and 212; Figure 2 also illustrates embodiments in which routers 203 205 controlled by the routing intelligence unit 200 are not coupled to SPALs. In Figure 3, a single routing intelligence unit 300 controls multiple edge routers 302 and 304, each of which is linked to exactly one ISP 306 and 308. In additional embodiments illustrated in Figure 4, different routing intelligence units 400 and 402, each connected to a set of local edge routers 404, 406, 408, and 410, may coordinate their decisions. In some embodiments of the invention, the routing intelligence units comprise processes running within one or more processors housed in the edge routers. Other configurations of routing intelligence units and edge routers will be apparent to those skilled in the art.

#### B. Architecture of Routing Intelligence Units

The routing intelligence units include a Decision Maker resource. At a high level, the objective of the Decision Maker is to improve the end-user, application level performance of prefixes whenever the differential in performance between the best route and the default BGP route is significant. This general objective has two aspects:

- One goal is to reach a steady state whereby prefixes are, most of the time, routed through the best available Service Provider Access Link (i.e., SPAL), that is, through the SPAL that is the best in terms of end-

to-end user performance for users belonging to the address space corresponding to that prefix. To achieve this goal, the Decision Maker will send a significant amount of updates to the router (over a tunable period of time) until steady state is reached. This desirable steady state results from a mix of customer-tunable criteria, which may include but are not limited to end-to-end user measurements, load on the links, and/or cost of the links.

- Current measurements of end-to-end user performance on the Internet show that fluctuations in performance are frequent. Indeed, the reasons for deterioration of performance of a prefix may include, but are not limited to the following:

The network conditions can vary along the path used by the packets that correspond to that prefix on their way to their destination.

Alternatively, the access link through which the prefix is routed can go down.

The Service Provider to which the prefix is routed can lose coverage for that prefix.

In such occurrences, the routing intelligence unit should detect the deterioration/failure, and quickly take action to alleviate its effect on the end-user.

In order to optimize application performance, the routing intelligence unit converts measurements on the performance of routes traversing the edge-routers into scores that rate the quality of the end-to-end user experience. This score depends on the application of interest, namely voice, video and HTTP web traffic. In some embodiments of the invention, by default, the routing intelligence unit attempts to optimize the performance of web applications, so its decisions are based on a score model for HTTP. However, in such embodiments, the customer has the choice between all of voice, video, and HTTP.

In order to avoid swamping routers with BGP updates, in some embodiments of the invention, the maximum rate of update permitted by the routing intelligence unit is offered as, for example, a control, such as a knob that

is set by the customer. The faster the rate of updates, the faster the system can react in the event of specific performance deteriorations or link failures.

However, the rate of updates should be low enough not to overwhelm the router. In some embodiments, the selected rate will depend on the customer's setting (e.g., the traffic pattern, link bandwidth, etc.); for example, faster rates are reserved to large enterprises where the number of covered prefixes is large. Even when the rate of updates is slow, in some embodiments of the invention, the most urgent updates are still scheduled first: this is performed by sorting the prefix update requests in a priority queue as a function of their urgency. The priority queue is then maintained in priority order. In some embodiments of the invention, the most urgent events (such as loss of coverage, or link failure) bypass this queue and are dealt with immediately.

In case interface statistics are available, the Decision Maker may directly use the corresponding information to function in an optimized way. For example, in some embodiments of the invention, the Decision Maker can use bandwidth information to make sure that a link of lower bandwidth is not swamped by too much traffic; in a similar manner, link utilization can be used to affect the rate of BGP updates sent to the router. Finally, the decision maker may use per-link cost information, as provided by the user to tailor its operation. For example, assume that the router is connected to the Internet through two links: Link 1 is a full T3, while Link 2 is a burstable T3, limited to 3 Mbit/sec. That is, whenever load exceeds the 3 Mbit/sec mark on Link 2, the user incurs a penalty cost. Combining information pertaining to per-link cost and utilization, the Decision Maker can attempt to minimize the instances in which load exceeds 3 Mbit/sec on Link 2, thus resulting in reduced costs to the user.

In some implementations, the Decision Maker may also use configurable preference weights to adjust link selection. The cost of carrying traffic may vary between links, or a user may for other reasons prefer the use of certain links. The Decision Maker can attempt to direct traffic away from some links and towards others by penalizing the measurements obtained on the less preferred links; conversely, if different links have comparable measured performance, traffic is directed away from the less preferred links.



Some embodiments of this invention can take into account more parameters, such as more information about SPALs and prefixes. However, despite the utility of such enhancements, the Decision Maker is designed to work well even when it relies on information provided by solely by the edge stats measurements.

In case the routing intelligence unit fails, the design is such that the edge router falls back to the routing that is specified in the BGP feed. The same Behavior takes place in case performance routes sent by the prefix scheduler Are filtered by the edge routers it controls.. Finally, in some embodiments of the invention, a flapping control algorithm is included in the design, avoiding the occurrence of undesirable excessive flapping of a prefix among the different access links.

A diagram showing the high-level architecture of Routing Intelligence Unit, and focused on its BGP settings is shown in Figure 5a. In the embodiments illustrated in Figure 5a, three BGP peering types may exist between a given Routing Intelligence Unit 500 and the external world: one to control the local edge router or routers 502 that this particular Routing Intelligence Unit 500 is optimizing, one to a Routing Infrastructure Exchange (RIX) 504, and one to every other Routing Intelligence Unit device with which it coordinates 506, as further described in U.S. Provisional Applications No. 60/241,450, filed October 17, 2000 and U.S. Provisional Application No. 60/275,206, filed March 12, 2001, U.S. Applications No. 09/903,441, filed July 10, 2001, U.S. Application No. 09/923,924, filed August 6, 2001, and U.S. Application No. 09/903,423, filed July 10, 2001, which are hereby incorporated by reference in their entirety. In the diagram shown in Figure 5a, the three external peering types are shown as the arrows at far left (to the Edge Routers 502 and to RIX 504) and far right 506. In order for BGP updates to be propagated to the appropriate devices, some devices are configured to be route reflectors, and others as route reflector clients. In embodiments illustrated in Figure 5a, the Edge Routers 502 are both route reflectors, and the peer BGP stacks are clients, as indicated by the labels "r" and "c". Similarly, in the peering between the BGP Process 506 and the BGP Stack, the BGP Process 506 is a route reflector, and the BGP Stack is a client. Note that the separation

between the BGP Process 506 and BGP Stack is not required in all embodiments. However, when they are separate, the use of route reflection allows the BGP Process 506 to behave as a normal BGP implementation (as described in The Big Book of Border Gateway Protocol RFCs referenced in the Background Of The Invention). Other configurations of the devices that may be used for propagation of BGP updates will be apparent to those skilled in the art.

### C. Coordination Between Routing Intelligence Units

Figure 5B schematically illustrates a configuration in which multiple routing intelligence units may coordinate via a back-channel to exchange routing information and set routing policy. Each Routing Intelligence Unit includes a Decision Maker 508 510 512, which in turn controls one or more routers 514 516 518 520 522. The routers 514 516 518 520 522 may in turn be coupled to one or more ISPs 524 526 528. Figure 5B also illustrates the back-channel 530, comprised of peerings between processes on Remote Coordination Processors (RCPs) 532 534 536; in some embodiments, these may be iBGP or eBGP peerings. Other implementations will be apparent to those skilled in the art. The back-channel 530, or mesh, may be used to communicate information on local path performance characteristics between Routing Intelligence Units, to increase the number of paths considered during optimization.

Such embodiments of the invention may employ BGP environments to support coordination between routers 514 516 518 520 522; alternatively, in some embodiments, this may be accomplished without BGP, by coupling the routers together, either physically or virtually. In embodiments of the invention utilizing BGP environments for coordination, the peerings on the back-channel 530 may be iBGP peerings.

In some embodiments of the invention, each of the Routing Intelligence Units sends its best local score to the others via the back-channel 530. In some such embodiments, local links are preferred over equivalent remote links. Additionally, in some such embodiments, a Routing Intelligence Unit does not send updates directly to remote routers. Rather, remote information is assessed by the local Routing Intelligence Unit prior to being forwarded to the associated router. In embodiments of the back-channel 530 utilizing BGP, techniques such

as route reflection and confederation may be used to scale the mesh. In one such embodiment, the coordination BGP processes may be arranged to match the original router BGP mesh as closely as possible, controlling each BGP router with a separate Routing Intelligence Unit. Other arrangements for the back-channel will be apparent to those skilled in the art.

In some embodiments of the invention, the routers under the control of the Decision Makers 508 510 512 are able to route between themselves by use of a single IP next-hop. For instance, in the example illustrated in Figure 5B, if a first router 514 forwards packets towards an established next-hop associated with a second router 518, then the packets will arrive at the second router 518.

In some embodiments, the Routing Intelligence Units coordinate by exchanging their best scores with one another. In some implementations, a Decision Maker 508 inside a Routing Intelligence Unit can elect to send an update on the back channel 530. In some such embodiments, this may occur whenever the Decision Maker 508 is also asserting to its routers 514 516. It may also occur when the Decision Maker 508 decides the current routes are correct. By exchanging information via the back channel 530, Decision Makers 508 510 512 may inform one another about local conditions. Additionally, if local scores change by a sufficient amount, this may be announced via the back-channel 530, even if the change in score doesn't affect local routing. In embodiments of the invention, the BGP processes used for coordination do not peer directly to the routers 514. Rather, they connect to the Decision Maker 508, and the Decision Maker 508 decides whether to pass on the update to the routers 514 516, as well as whether to modify it.

In some embodiments of the invention, the BGP process for coordination is configured so that the Decision Maker 508 is a route reflector client of the other Decision Makers 510 512. The Decision Maker 508 is also a route reflector client of the edge routers it controls 514 516. Thus, in such embodiments, the Decision Makers 508 510 512 do not simply transmit information in either direction without consideration; rather, these BGP processes are separate data channels.

In embodiments of coordination implemented with BGP, a scalar performance score exchanged between Routing Intelligence Units may be

translated to units of Local Preference, where some implementations of Local Preference use 8 bits and others use 16 bits. Using Local Preference ensures that the new BGP mesh 530, or back-channel, will automatically select and propagate the best score. Other embodiments of the invention implemented with BGP may transfer scalar performance scores encoded within the community attribute, the extended communities attribute, the multi-exit discriminator attribute, or some combination of all of the above.

Embodiments of the invention also include procedures for a Decision Maker 508 to decide whether to use a prefix which arrives via coordination with the other decision makers 510 512. Some implementations avoid use of such remote routes unless they are distinctly attractive. Thus, in such embodiments, given a choice between comparable local and remote routes (wherein 'comparable' may mean within a winner-set width), the local route is always used. Other implementations may include:

- a static penalty applied to all remote announcements
- a static penalty per remote Decision Maker
- a static penalty per remote SPAL
- dynamic penalties per remote Decision Maker

In the case of dynamic penalties per Decision Maker, it is possible to have one Decision Maker 508 probe all others 510 512 actively, and use the measure of distance between Routing Intelligence Units as a dynamic penalty. Other methodologies for implementing dynamic penalties will be apparent to those skilled in the art.

#### D. Queuing Architecture

A diagram showing the high level mechanics of the decision maker prefix scheduler is shown in Figure 6. As illustrated in Figure 6, two threads essentially drive the operation of the scheduler. The first thread polls the database for changes in terms of per-SPAL performance, load, or coverage, and decides on which prefix updates to insert in a Priority Queue that holds prefix update requests. The second thread takes items out of the queue in a rate-controlled fashion, and converts the corresponding update requests into an appropriate set of UPDATES that it sends to the local routers, and an

appropriate set of UPDATES that it sends to the back channel for communication to other Routing Intelligence Units.

In the following, we describe each thread separately. In the description, we will refer to tables in the database, and to fields within these tables. The contents of this database are also explicated in U.S. Provisional Applications No. 60/241,450, filed October 17, 2000 and U.S. Provisional Application No. 60/275,206, filed March 12, 2001, and U.S. Applications No. 09/903,441, filed July 10, 2001, U.S. Application No. 09/923,924, filed August 6, 2001, and U.S. Application No. 09/903,423, filed July 10, 2001, which are hereby incorporated by reference in their entirety.

#### Thread 1

This first thread 600 polls the database for changes in terms of per-SPAL performance, load, or coverage, and decides on which prefix updates to insert in a Priority Queue that holds prefix update requests.

In some embodiments of the invention, such changes are checked for in 2 passes. The first pass looks for group level changes, wherein a group comprises an arbitrary collection of prefixes. Groups are also described in U.S. Provisional Applications No. 60/241,450, filed October 17, 2000 and U.S. Provisional Application No. 60/275,206, filed March 12, 2001, and U.S. Applications No. 09/903,441, filed July 10, 2001, U.S. Application No. 09/923,924, filed August 6, 2001, and U.S. Application No. 09/903,423, filed July 10, 2001, which are hereby incorporated by reference in their entirety. In case a significant change in performance for a group is noticed, the group is unpacked into its individual prefixes; the corresponding prefixes are checked and considered for insertion in the priority queue. The second pass captures prefixes for which there are no group-level performance changes.

An update request for a prefix can be made in a number of different circumstances. Non-limiting examples of such circumstances include any one or more of the following:

- 1) In case a significant change in its performance score is witnessed on at least one of its local SPALs.

- 2) In case a significant change in its performance score is witnessed on a foreign SPAL (that is, a SPAL that is controlled by a different Routing Intelligence Unit box in a coordinated system).
- 3) In case any of the local SPALs becomes invalid.
- 5 4) In case an update pertaining to this prefix was received from the router.
- 5) A peering with either a local or a remote router goes down, for instance, during the router's maintenance windows.
- 6) At the user's request.

10 Note that measurements reside at the group level; hence, Check 1 can be done in the first pass. On the other hand, all of Checks 2, 3, and 4 are prefix-specific and may be performed in Pass 2: indeed, foreign performance updates are transferred through the back channel in BGP messages, and hence correspond to particular prefixes. Also, SPALs may become invalid for some, and not  
 15 necessary all prefixes in a group. Finally, updates from the router relate to the change of winner SPALs for some prefixes, or to the withdrawal of other prefixes. (In fact, any information that is transferred by BGP relates to prefixes.)

#### Pass 1:

20 In some embodiments of the invention, in the first pass, an asynchronous thread goes through all groups in the GROUP\_SPAL table, checking whether the NEW\_DATA bit is set. This bit is set by the measurement listener in case a new measurement from a /32 resulted in an update of delay, jitter, and loss in the database. Delay, jitter, and loss, also denoted as  $d$ ,  $v$ , and  $p$ , are used to compute an application-specific score, denoted by  $m$ . The scalar  $m$  is used to  
 25 rate application-specific performance; MOS stands for "Mean Opinion Score", and represents the synthetic application-specific performance. In embodiments of the invention, MOS may be multiplied by a degradation factor that is function of link utilization, resulting in  $m$ . (That is, the larger the utilization of a given SPAL, the larger the degradation factor, and the lower the resulting  $m$ )

30 In embodiments of the invention, users of the device may also configure penalty factors per SPAL. Non-limiting examples of the uses of such penalty features include handicapping some links relative to others, to achieving cost

control, or accomplishing other policy objectives. As a non-limiting example, Provider X may charge substantially more per unit of bandwidth than Provider Y. In such a situation, the penalty feature allows the user to apply an **m** penalty to SPAL X. This will cause Provider Y to receive more traffic, except for those  
 5 prefixes in which the performance of Provider X is substantially better. One implementation of this embodiment is to subtract the penalty for the appropriate SPAL after **m** is computed. Other implementations of the penalty feature will be apparent to those skilled in the art.

Even when NEW\_DATA is set, the variation in d, v, and p can be small  
 10 enough so that the change in the resulting scalar **m** is insignificant. Hence, in some embodiments of the invention, the prefix is only considered for insertion in the queue in case the change in **m** is significant enough. The corresponding pseudo-code is shown below.

```

for each group
15  {
        // First pass: only consider groups for
        which there is a change in the group pref data
        compute_winner_set = 0;

        for each spal (<> other)
20      {
            // check whether there is new data for
            this group
            if (new_data(group, spal)==1)
25          {
                compute m (spal, d, v, p, spal-
                penalty), store in local memory
                new_data(group, spal) = 0;
                if (significant change in m)
30          {
                        store m (spal, d, v, p)
                        in group_spal

                        compute_winner_set = 1;
                        break;
  
```

```

        }
    }
}

5         if (compute_winner_set)
            for each prefix
                schedule_prefix(prefix) // see
below
    }
10
    In some embodiments of the invention, rolling averages are used to
update measurements of delay, jitter, and loss, i.e.,
        d = alpha*d + (1 - alpha)*dnew
        v = beta*v + (1 - beta)*vnew
        p = gamma*p + (1 - gamma)*pnew,
15
where dnew, vnew, pnew represent the new delay, jitter, and loss
measurements. Algorithms for calculating MOS for HTTP (1.0 and 1.1) and for
voice and video are also presented in U.S. Provisional Applications No.
60/241,450, filed October 17, 2000 and U.S. Provisional Application No.
60/275,206, filed March 12, 2001, and U.S. Applications No. 09/903,441, filed
20 July 10, 2001, U.S. Application No. 09/923,924, filed August 6, 2001, and U.S.
Application No. 09/903,423, filed July 10, 2001. Values used for the models
employed by these algorithms in embodiments of the invention are presented in
an XML format below. Note that since MOS is computed per group, a selection
from the sets of the following parameters may be made to allow different
25 optimization goals for each group.

<module> <engine slot="1"> <application model="http1.0" [alpha="0.9"
beta="0.9" gamma="0.9" theta="1.18" phi="0.13" omega="0.15" psi="0.25"]
/>
30 </engine> </module>

<module> <engine slot="1"> <application model="http1.1" [alpha="0.9"
beta="0.9" gamma="0.9" theta="1.3" phi="0.31" omega="0.41" psi="1.0"] />
</engine> </module>

```



```
<module> <engine slot="1"> <application model="voice"
[alpha="0.9" beta="0.9" gamma="0.9" theta="1.5" phi="6.0" omega="23.0"
psi="0.0"] /> </engine>
```

5 </module>

```
<module> <engine slot="1"> <application model="video" [alpha="0.9"
beta="0.9" gamma="0.9" theta="1.0" phi="4.0" omega="69.0" psi="0.0"] />
</engine> </module>
```

10 The values presented above are given as examples only. Many different models for deriving MOS scores for different applications will be apparent to those skilled in the art.

### Pass 2

In some embodiments of the invention, in the second pass, an asynchronous thread goes through all prefixes in the PREFIX table. In some such embodiments, for each prefix, Checks 2, 3, and 4 are made: NEW\_INCOMING\_BID in the PREFIX table indicates that a new bid was received from the coordination back channel; NEW\_INVALID in the PREFIX\_SPAL table indicates, for a particular (Prefix P, SPAL x) pair a loss of coverage for Prefix P over SPAL x. NEW\_NATURAL\_DATA indicates the receipt by Routing Intelligence Unit of an update message from a router, notifying it of a change in its natural BGP winner. In fact, the Decision Maker only asserts a performance route in case it is not the same as the natural BGP route; hence, it can potentially receive updates concerning the natural BGP winners of given prefixes from routers to which it has asserted no performance route for those prefixes. (The advantage of such an implementation is that when no performance route is sent to a router, the routing intelligence unit will get routing updates from that router. In contrast, if performances route were asserted regardless of whether they agree with the natural BGP choice, the Routing Intelligence Unit would never receive an update from the router pertaining to changes in the natural BGP winner for the different prefixes. If Routing Intelligence Unit were to assert performance

routes regarding a given prefix P to all routers irrespective of the current BGP winner for that prefix, it will never receive an update from the router pertaining to changes in the natural BGP winner for Prefix P. Indeed, the performance route would always be the winner, so the router would assume there is nothing to talk about.)

The following example illustrates the usefulness of the NEW\_NATURAL\_DATA flag: Assume that the Decision Maker controls 3 routers, each of which controls its individual SPAL. Assume that the Decision Maker has just determined that Prefix P will move to SPAL 1. Assume that Prefix P believes that the natural BGP route for Prefix P as saved by Router 1 is SPAL 1, the same as its current performance assertion. The Decision Maker's logical operation is to withdraw Prefix P's last performance route (say SPAL 3). However, it turned out that this BGP natural route has, in fact changed to SPAL 2; indeed, this could have happened during the previous assertion of a performance route for Prefix P (since, in this case, as mentioned above, the Decision Maker receives no updates for Prefix P from the router, despite potential changes in Prefix P's natural BGP winner). As a result of this discrepancy, all traffic pertaining to Prefix P will be routed through SPAL 2, the current natural BGP winner for Prefix P, which is not the desired behavior.

This is the primary reason for NEW\_NATURAL\_DATA: as such an event occurs, the router sends an update back to the Decision Maker, communicating to it the change in natural route. The incoming BGP messages from the local routers are processed by a process referred to as the Peer Manager. The Peer Manager sees the change in natural BGP route and sets the NEW\_NATURAL\_DATA flag to 1; consequently, the prefix is considered for re-scheduling during this pass, in Thread 1, as described above. Note that in case of changes in the natural BGP route for a given prefix, the Decision Maker will need two passes through the Priority Queue before the prefix is routed through its appropriate performance route.

Finally, the ACCEPTING\_DATA bit in the prefix table is checked. ACCEPTING\_DATA is set to 0 by the peer manager to notify the decision maker not to assert performance routes for this prefix. This would primarily

occur in case the prefix is withdrawn from the BGP tables in all local routers. In this case, in the ROUTER\_PREFIX\_SPAL table, the peer manager would have set the ANNOUNCED bits for that prefix on all SPALs to zero. Clearly, a prefix is only considered for insertion in the queue in case

```

5  ACCEPTING_DATA is set to 1.

    for each prefix
    {
        //Checks 2 and 4: scan the prefix_group table
10      get new_bid, new_natural, and accepting_data from
        prefix_group
        if (new_bid) || (new_natural)
        {
            if (accepting_data)
15          {
                schedule_prefix(prefix) // see
            below
            }
        }
        //Check 3: scan the prefix_spal table
20      get new_invalid, from prefix_spal
        if (new_invalid)
        {
            schedule_prefix(prefix)
        }
25  }

```

Note that asserting a performance route about a prefix that does not exist in any of the routers' BGP tables could be problematic, depending on the surrounding network environment. If the set of controlled routers do not emit routes to any other BGP routers, then it is acceptable to generate new prefixes. But if any propagation is possible, there is a danger of generating an attractor for some traffic.

Specifically, if the new route is the most specific route known for some addresses, then any traffic to those addresses will tend to forward from uncontrolled routers towards the controlled routers. This can be very disruptive, since such routing decisions could be very far from optimal.

The mechanism can cope with this in a number of ways:

- Prevent any use of a prefix unknown to BGP. This is achieved using the ACCEPTING\_DATA check included in some embodiments of the invention.
- 5       • Permit all such use, in a context where new routes cannot propagate
- Permit such use, but mark any new prefix with the well-known community value no-advertise to prevent propagation
- 10       • Permit such use, but configure the routers to prevent any further propagation (in some embodiments, by filtering such prefixes)

Deciding to Insert a Prefix Update Request in the Priority Queue:

The *schedule\_prefix* Function

Once a prefix P makes it through the checks imposed in either Pass 1 or Pass 2, it is considered for insertion into the prefix update priority queue.

15       *schedule\_prefix* includes the related functionality, described below:

- First of all, a winner set of SPALs is re-computed for P; this set includes SPALs for which the performance is close to maximal.
- After the winner set W is computed for P, the decision maker determines whether the current route for P is included in W.
- 20       • In case of a coordinated Routing Intelligence Unit system, in some embodiments of the invention, the back channel is sent updates pertaining to Prefix P even if the local prefix update request is dropped. For example, the performance on local links could have changed dramatically since the last time a bid was sent to the back channel for this prefix; in the event of such an occurrence, an updated bid is sent to the back channel (through the BGP peering set up for this purpose).
- 25       • In case the current route is not part of the newly computed winner set, it is clear that Prefix P is not routed optimally. Before going ahead and inserting an update request for Prefix P in the queue, the Routing Intelligence Unit performs a check of the flapping history for Prefix P.
- 30

In case this check shows that Prefix P has an excessive tendency to flap, no prefix update request is inserted in the queue.

- In some embodiments of the invention, before the prefix is inserted in the queue, a SPAL is chosen at random from the winner set. In case the winner set includes a remote SPAL controlled by a coordinated Routing Intelligence Unit as well as a local SPAL, the local SPAL is always preferred. Also, in some embodiments of the invention, the randomness may be tweaked according to factors pertaining to any one or more of the following: link bandwidth, link cost, and traffic load for a given prefix. Finally, the state in the database is updated, and the element is inserted in the Priority Queue. The rank of the prefix update in the priority queue is determined by computing the potential percent improvement obtained from moving the prefix from its current route to the pending winner route.

At the outset, a winner set of SPALs is re-computed for P; this set includes SPALs for which the performance is close to maximal. In some embodiments of the invention, invalid SPALs are excluded from the winner set computation. Bids from remote SPALs under the control of coordinated Routing Intelligence Units may, in embodiments, be included in the winner set computation. Since the bids corresponding to such remote routes are filtered through BGP, they are in units which are compatible with iBGP's LocalPref, which in some implementations is limited to 0-255. Therefore one possible implementation is to multiply  $m$  by 255. The converted quantity is referred to as MSLP. For consistency, the  $m$  values computed for local SPALs are also converted to local\_pref units. The new winner is then determined to be the set of all SPALs for which MSLP is larger than  $MSLP_{max} - \text{winner-set-threshold}$ , where  $MSLP_{max}$  represents the maximum MSLP for that prefix across all available SPALs, and winner-set-threshold represents a customer-tunable threshold specified in LocalPref units. The related pseudo-code is shown below.

```
for each spal (<> other)
{
```

```

        get invalid bit from prefix_spal
        if (invalid)
        {
            mark spal as invalid, not to be used in
5 winner_set computation
            continue
        }
        convert m (spal) to MSLP
        Store MSLP in prefix_spal table
10 }
    for spal=other
    {
        get MSLP_other = other_bid in prefix_group table
    }
15 compute winner_set(prefix) // considers winners among all
    valid spals and other_bid

```

After the winner set W is computed for P, the decision maker determines whether the current route for P is included in W. Indeed, in such a case, the performance of that prefix can't be improved much further, so no prefix update request needs to be inserted in the queue.

Even though an update request for a given prefix is ignored, the Decision Maker may still send an update to the back channel in certain embodiments. For example, even though the current route for Prefix P is still part of the winner set, performance degradation could have affected all SPALs at once, in which case the bid that was previously sent to the back channel for Prefix P is probably inaccurate. In some embodiments, one may solve this problem by implementing the following: the last bid for a given prefix is saved as MY\_BID in the PREFIX table; a low and high threshold are then computed using two user-configurable parameters, bid-threshold-low and bid-threshold-high. In case of a significant difference between the MSLP score on the current route and the last score sent to the back channel for that prefix (i.e., MY\_BID) is witnessed (that is, if the new score falls below  $(1 - \text{bid-threshold-low}) * 100\%$  or jumps to a value that is larger than  $(1 + \text{bid-}$

threshold-high)\*100% of MY\_BID), a BGP message is sent to the back channel, carrying the new bid for Prefix P to remote coordinated Routing Intelligence Units. Pseudo-code illustrating the functionality described here is shown below.

```

5
//First, detect non-communicated withdrawal of a prefix
if winner_set only comprises remote link
{
    for all local routers
10         if performance route exists for that
        (prefix, router) pair in the ROUTER_PREFIX_SPAL table
            send urgent withdrawal of this
            route to edge router
            continue
15 }
get current_winner(prefix) and pending_winner(prefix) from
prefix_spal table

if (pending_winner!=current_winner)
20 {
    if (current_winner in winner_set)
    {
        update pending_winner = current_winner in
        database
25         continue
    }
    if (current_winner not in
winner_set)&&(pending_winner in winner_set)
    {
30         continue
    }

}

35
if (current_winner==pending_winner)

```

```

{
    if (new_natural)
    {
        for all routers
5         {
            current_route_per_router =
            SPAL(prefix, router, type = natural, state = latest_ON)
            if (current_route_per_router
            exists) && (current_route_per_router != current_winner)
10         {
            special_route =
            current_route_per_router
            set local
            special_route_flag = 1;
15         break;
        }
    }
    else
    {
20         current_route = current_winner
    }
    if (current_route in
    winner_set) || (special_route == current_winner)
25     {
        get bid_low_threshold and
        bid_high_threshold from prefix_group table
        if ((MSLP(prefix, current_spal) <
        bid_low_threshold) || (MSLP(prefix, current_spal)
30     bid_high_threshold))
        {
            compute bid_low_threshold and
            bid_high_threshold from MSLP(prefix)
            store bid_low_threshold and
35     bid_high_threshold in prefix_group

```



```

                                form UPDATE to send to backchannel
SBGP
                                }
                                continue
5      }
  }

```

At this point, it is clear that Prefix P is not routed optimally. In some embodiments of the invention, before proceeding with sending the update request to the edge router, the Routing Intelligence Unit performs a check of the flapping history for Prefix P. An algorithm whose operation is very close to the flapping detection algorithm in BGP monitors the flapping history of a prefix. The algorithm can be controlled by, in one embodiment, three user-controlled parameters `flap_weight`, `flap_low`, and `flap_high` and works as follows: the tendency of a prefix to flap is monitored by a variable denoted `FORGIVING_MODE` that resides in the `PREFIX` table. `FORGIVING_MODE` and other flapping parameters are updated in Thread 2 right before a performance route pertaining to Prefix P is asserted to the local routers. In case `FORGIVING_MODE` is set to 1, the tendency for Prefix P to flap is considered excessive, and the prefix update request is ignored. Conversely, in case `FORGIVING_MODE` is set to 0, Prefix P has no abnormal tendency to flap, so it is safe to consider its update request.

```

get flapping state for prefix from prefix_group table
25  if (excessive flapping)
    {
        continue
    }

```

30 If a prefix survives to this point in Thread 1, it will deterministically be inserted in the queue. Hence, all bits that were checked should be reset at this point so that some other pass on the prefixes does not reconsider and reschedule the prefix update request. For example, in case the prefix belongs to a group for

which there was a significant change in **m**, the prefix will be considered for insertion in the queue in Pass 1, and should not be reconsidered in Pass2.

```

//reset prefix level bits, if necessary
5  for each spal (<> other)
    {
        get new_invalid bit from prefix_spal
        if (new_invalid)
            reset new_invalid to 0 in prefix_spal
10  }
    get new_bid and new_natural bits from prefix_group
    if (new_bid)
        reset new_bid to 0 in prefix_group
    if (new_natural)
15  reset new_natural to 0 in prefix_group

```

In some embodiments of the invention, before the prefix is inserted in the queue, a SPAL is chosen at random from the winner set. This way, traffic is spread across more than one SPAL, hence achieving some level of load

20 balancing. In order to achieve some set of desirable policies, randomness can be tweaked in order to favor some SPALs and disregard others. For example, in some embodiments, in case the winner set includes a remote SPAL controlled by a coordinated Routing Intelligence Unit as well as a local SPAL, the local SPAL is always preferred. In other words, a remote SPAL is only the winner in

25 case it is the only available SPAL in the winner set. Also, depending on the weight of a prefix and the observed load on different links, one can tweak the probabilities in such a way that the prefix is routed through a SPAL that fits it best. (This feature corresponds to the "Saturation Avoidance Factor" – SAF, described later in this document) After a winner is selected,

30 PENDING\_WINNER in PREFIX\_SPAL is updated to reflect the new potential winner. Finally, the element is inserted in the Priority Queue. In some embodiments, the rank of the prefix update in the priority queue is determined by computing the percent improvement; that is, the percent improvement obtained from moving the prefix from its current route to the pending winner

route. That is,  $\text{percent\_improvement} = [\text{score}(\text{pending\_winner}) - \text{Score}(\text{current\_route})] / \text{Score}(\text{current\_route})$ . The special-spal-flag is part of the data structure for the update, as it will be used in the determination of which messages to send to the local routers.

```

5
    if ((winner_set_size>1) and (other in winner_set))
        remove other from winner_set
    select spal from winner_set at random
    update PENDING_WINNER in PREFIX_SPAL table
10    compute percent_improvement for prefix
    insert prefix in prefix update queue

```

### Thread 2

In this thread 702, elements are taken out of the queue in a rate-controlled manner. In some embodiments of the invention, this rate is specified by the customer. The update rate is often referred to as the token rate. Tokens are given at regular intervals, according to the update rate. Each time a token appears, the head of the queue is taken out of the queue, and considered for potential update. In case the database shows that more recent passes in Thread 1 have canceled the update request, it is dropped *without losing the corresponding*

20 *token*; the next update request is then taken out from the head of the queue; this procedure is performed until either the queue empties, or a valid request is obtained. In some embodiments of the invention, when an update request that corresponds to Prefix P is determined to be current (thus, valid), one or more of the following tasks are performed:

25     The flapping state is updated for Prefix P.

      The database is updated to reflect the new actual winner; more specifically, the pending winner, chosen before inserting the prefix update request at the end of the first thread now becomes the current winner.

30     The database is checked to determine the current state of each of the individual routers. Accordingly, individual UPDATES are formed and sent to each of the routers. For example, no performance route is sent to an edge router

in case the BGP winner for Prefix P, according to that router is found to be the same.

An UPDATE is sent to the back channel, describing the new local winner.

5 Finally, the database is updated to keep track of the messages that were sent to each of the routers, as well as the expected resulting state of these routers.

In this thread 702, elements are just taken out from the queue in a rate-controlled manner, according to an update rate that may be set by the customer. The update rate is often referred to as the token rate: indeed, tokens are given at regular intervals, according to the update rate. Each time a token appears, the  
10 head of the queue is taken out, and considered for potential update.

Assume that the update request concerns Prefix P. The PREFIX\_SPAL table is checked to obtain the PENDING\_WINNER and CURRENT\_WINNER for Prefix P.

In case PENDING\_WINNER and CURRENT\_WINNER correspond to the  
15 same SPAL, this is an indication that a more recent pass in Thread 1 has canceled the update request; in this case, the update request is dropped, *without losing the corresponding token*; the next token request is then polled from the head of the queue; this procedure is performed until either the queue empties, or a valid request, for which PENDING\_WINNER and CURRENT\_WINNER are  
20 different, is obtained.

Having different pending and current winners reflects a valid update request. In this case, the Decision Maker should assert the winning route for Prefix P. When a prefix update request is considered still valid, it is implemented. In the process, a series of tasks are performed. First, the flapping  
25 state is updated for Prefix P. In some embodiments of the invention, the tendency of a prefix to flap is monitored by a variable denoted INTERCHANGE\_RATE that resides in the PREFIX table. The `flap_weight` parameter dictates the dynamics of INTERCHANGE\_RATE; more specifically, at this point in the algorithm thread, INTERCHANGE\_RATE is updated using  
30 the last value of INTERCHANGE\_RATE, as stored in the table, LAST\_ICR\_TIME, also stored in the PREFIX table, and `flap_weight`. In case the new computed INTERCHANGE\_RATE is below `flap_low`, Routing Intelligence Unit considers the tendency for that prefix to flap to be low. On the

other hand, when INTERCHANGE\_RATE exceeds `flap_high`, the Routing Intelligence Unit considers the tendency for that prefix to flap to be high. That is, the algorithm functions in the following fashion:

- 5       • In case FORGIVING\_MODE (also in the PREFIX table) is set to 0, and INTERCHANGE\_RATE exceeds `flap_high`, FORGIVING\_MODE is set to 1.
- In case FORGIVING\_MODE is set to 1, but INTERCHANGE\_RATE drops below `flap_low`, FORGIVING\_MODE is set to 0 again, and the prefix update request survives this check.
- 10     • In case FORGIVING\_MODE is set to 1 and INTERCHANGE\_RATE is larger than `flap_low`, or FORGIVING\_MODE is set to 0, and INTERCHANGE\_RATE is below `flap_high`, FORGIVING\_MODE does not change.

Note that the method presented above is only one technique for controlling flapping; others will be apparent to those skilled in the art.

15       In some embodiments of the invention, the two parameters `flap_low`, and `flap_high` are separated by an amount to avoid hysteresis between the two values. Then, the Decision Maker updates the PREFIX\_SPAL table to reflect this change; more specifically, CURRENT\_WINNER is moved to  
20       PENDING\_WINNER in the table. At this time, the ROUTER\_PREFIX\_SPAL table is queried to capture the current state of each router in regards to Prefix P. Accordingly, different UPDATES are formed and sent to each of the routers.

      In some embodiments of the invention, the Decision Maker only asserts a performance route in case it is not the same as the natural BGP route; indeed,  
25       if Routing Intelligence Unit were to assert performance routes regarding a given prefix P to all routers irrespectively of the current BGP winner for that prefix, it will never receive an update from the router pertaining to changes in the natural BGP winner for Prefix P. (Indeed, the performance route would always be the winner, so the router would assume there is nothing to talk about.)

30       Also, an UPDATE is sent to the back channel, describing to other Routing Intelligence Units in a coordinated system the new local winner. Finally, the database is updated to keep track of the messages that were sent to each of the routers, as well as the expected resulting state of these routers.

Prior to forming the UPDATES, the database is updated as to include the new flap parameters and prefix-SPAL information (i.e., the new current SPAL for that prefix). The BGP update sent to an edge router may be filtered out by policy on the router. However, assuming the update is permissible, it may be made to win in the router's BGP comparison process. One implementation is to have the edge router to apply a high Weight value to the incoming update. (Weight is a common BGP knob, supported in most major implementations of the protocol, but it is not in the original protocol specification) This technique constrains the update so that it gains an advantage only on the router or routers to which the update is directly sent; this is desirable if some other routers are not controlled by a device such as the one described here. It is also possible to send the update with normal BGP attributes which make the route attractive, such as a high LocalPref value.

```

15  if (local_token available)
    {
        get prefix at the head of the local update queue
        updatePrefixSpal(prefix, spal)
20    updateFlapStats(prefix)
        compute bid_low_threshold and bid_high_threshold
        from MSLP(prefix)
        store bid_low_threshold and bid_high_threshold in
        prefix_group
25    form UPDATE to send to local SBGP
        form UPDATE to send to backchannel SBGP
    }

```

#### E. Technical Considerations

##### Queue Size

30 In some embodiments of the invention, a maximum queue size is to be chosen by the customer. In some embodiments, a small queue size may be chosen, so the maximum delay involved between the time instant a prefix update request is queued and the time instant it is considered by the second

thread as a potential BGP update is small. For example, in case the token rate corresponding to a given link is 10 tokens per second, and we choose not to exceed a 2 second queuing delay, the queue should be able to accommodate 20 prefix update requests. Note that this method is simple, and only requires the knowledge of the token rate and the maximum acceptable delay.

#### Maximum Rate of Prefix Updates

It is desirable for the Routing Intelligence Unit to remain conservative in the rate of updates it communicates to the edge-router. This is the function of the token rate, which acts as a brake to the whole system. In some embodiments of the invention, the responsibility for setting the token rate is transferred to the customer, who selects a token rate that best fits her bandwidth and traffic pattern.

#### F. Feedback from the Listener BGP

The feedback from the listener BGP is valuable as it describes the actual current state of the local edge routers. Accordingly, in some embodiments of the invention, a separate routing intelligence unit thread modifies the content of the database according to the state it gets from the router(s). The Routing Intelligence Unit can operate more subtly in case it is a *perfect listener*; we consider the Routing Intelligence Unit to be a perfect listener if it has knowledge of the individual BGP feeds from each individual SPAL. That is, in case the Routing Intelligence Unit is connected to three access links, each connecting to a separate provider, the Routing Intelligence Unit is a perfect listener if it has access to each of the three feeds handed by each of these providers.

Configuring Routing Intelligence Unit as a Perfect Listener is desirable, as it allows the support of private peerings. For example, unless Routing Intelligence Unit is configured as a Perfect listener, when Routing Intelligence Unit hears about a prefix, it can't assume that coverage exists for that prefix across all SPALs. Considering the scenario described above, a prefix that the Routing Intelligence Unit learns about could be covered by any of the three SPALs the router is connected to. For example, assume that only SPAL 1 has coverage for a given prefix P; in case the Routing Intelligence Unit asserts a

performance route for that prefix across SPAL 2, there is no guarantee that the traffic pertaining to that prefix will be transited by the Service Provider to which SPAL 2 is connected (which we denote Provider 2). In case Provider 2 actually has a private peering with Provider X that obeys to some pre-specified contract, Provider X could well monitor the traffic from Provider 2, and filter all packets that do not conform to that contract. In case this contract namely specifies that Provider X will only provide transit to customers residing on Provider X's network, then the traffic pertaining to Prefix P will be dropped. If Routing Intelligence Unit were a Perfect Listener, it would only assert performance routes for prefixes across SPALs that are determined to have coverage for these prefixes. This behavior may be referred to as "extremely polite".

In some embodiments, the Routing Intelligence Unit is capable of avoiding the "Rocking the boat" problem, which stems from unwanted propagation of prefixes which did not already exist in BGP. The Routing Intelligence Unit can operate in "impolite" mode, where any prefixes may be used, or in "polite" mode, where only those prefixes which were previously present in BGP can be used. An ANNOUNCED bit resides in the ROUTER\_PREFIX\_SPAL table, and is set by the Peer Manager in case the Routing Intelligence Unit hears about a prefix from any of the Routers. This bit allows use of "polite" mode by the following procedure: in case the ANNOUNCED bit is set to 0 for all (router, SPAL) combinations in the ROUTER\_PREFIX\_SPAL table, then ACCEPTING\_DATA is set to 0 in the PREFIX table.

25

#### G. Urgent Events

In case a catastrophic event occurs, such as a link going down, some embodiments of the invention send urgent BGP updates to the router. These urgent updates have priority over the entire algorithm described above. For example, in case a SPAL has lost coverage for a prefix, an urgent BGP message should be sent to the router, requesting to move the prefix to other SPALs. A

30



list of urgent events upon which such actions may be taken, and a description of the algorithms pertaining to these actions, are described below.

#### Algorithm for the Detection of an Invalid SPAL

In some embodiments of the invention, a specific (Prefix P, SPAL x)  
5 pair is invalidated in case there are reasons to believe that SPAL x no longer provides coverage to Prefix P. One possible implementation is described as follows. Measurements corresponding to a (Prefix, SPAL) pair are assumed to arrive to the Decision Maker at something close to a predictable rate. A background thread that is independent from Threads 1 and 2 computes this  
10 update rate, and stores a time of last update, the LAST\_UPDATE\_TIME. Another background thread verifies that LAST\_ICR\_TIME is reasonable given UPDATE\_RATE. For example, assuming that measurements come in following a Poisson distribution, it is easy to verify whether LAST\_ICR\_TIME exceeds a fixed percentile of the inter-arrival interval. As LAST\_UPDATE\_TIME  
15 increases, the Decision Maker becomes more and more worried about the validity of the path. In the current design, there are two thresholds: at the first threshold, the NEW\_INVALID and INVALID flags are set in the PREFIX\_SPAL table. As described in Thread 1 above, setting the NEW\_INVALID flag for a (Prefix P, SPAL x) pair will prevent any new update  
20 requests for Prefix P to be routed through SPAL x. At this stage, no other action is taken. At the second threshold, the Decision Maker becomes “very concerned” about routing Prefix P through SPAL x; hence, an urgent check is made to see whether Prefix P is currently routed through SPAL x, in which case an urgent UPDATE is created (that is, an UPDATE that bypasses the entire  
25 queue system) in order to route Prefix through a different SPAL.

#### H. Saturation Avoidance Factor

Some embodiments of the invention support a Saturation Avoidance Factor, which measures the effect of a prefix on other prefixes. In some  
embodiments of the invention, the “Saturation Avoidance Factor” (SAF)  
30 pertaining to a given prefix may be taken into account when prefixes are sorted in the Priority Queue. This SAF measures the effect of a prefix on other

prefixes. That is, if, upon scheduling a prefix on a given link, its effect on the other prefixes already scheduled on that link is high (i.e., this effectively means that the aggregate load for this prefix is large), its SAF should be low. The lower the SAF of a prefix, the lower its place in the Priority Queue. This way, the algorithm will always favor low load prefixes rather than high load prefixes. Note that in some embodiments, the SAF is not directly proportional to load. For example, a prefix that has a load equal to  $0.75C$  has a different SAF whether it is considered to be scheduled on an empty link or on a link which utilization has already reached 75%. In the later case, the SAF should be as low as possible, since scheduling the prefix on the link would result in a link overflow.

At times, the token rate may be slower than the responded feedback. In case link utilization information is obtained through interface-stats, the token rate may be slower than the rate at which utilization information comes in. Also, the token rate may be slower than the rate at which edge-stats measurements come in.

Additionally, in some embodiments, each prefix is considered at a time. That is, PQServiceRate is small enough so that no more than one token is handed at a time. For example, denoting by  $T$  the token rate obtained from the above considerations, PQServiceRate is equal to  $1/T$ . If more than one token were handed at one time, two large prefixes could be scheduled on the same link, just as in the example above, potentially leading to bad performance.

In some embodiments of the invention, the SAF is a per-prefix, per-SPAL quantity. For example, assume that a prefix carries with it a load of 75% the capacity of all SPALs. If we have a choice between two SPALs, SPAL 1 and SPAL 2, SPAL 1 already carrying a load of 50 %, the other having a load of 0%. In this case, moving Prefix  $p$  to SPAL 1 will result in bad performance not only for itself, but also for all other prefixes already routed through SPAL 1. In this case, the SAF is close to 0, even if performance data across SPAL 1 seems to indicate otherwise. On the other hand, the SAF of moving Prefix  $p$  to SPAL 2 is, by contrast, very good, since the total load on the link will remain around 75% of total capacity, so delays will remain low. If, instead of carrying a load of 75% capacity, Prefix  $p$  carried a load of 10% capacity, the results would have

been different, and the SAF of Prefix p across SPALs 1 and 2 would have been close. In some embodiments of the invention, without knowing the load of a link, we can still measure the effect of moving a given prefix to a given SPAL through RTT measurements. That is, instead of measuring the load directly, we  
5 measure the end result, that is the amount by which performance of prefixes across a link worsens as a result of moving a prefix to it.

#### Modifying the Schema for the Support of SAF

In order to support SAF, the schema may be include a load field in the SPAL table, and an SAF field in the PREFIX\_SPAL table. In some  
10 embodiments, the SAF field is a per-prefix, per-SPAL information.

#### I. Available Bandwidth

Edge-stats measurements may include measurements of delay, jitter, and loss; using these measurements, an application-specific performance score may  
15 be obtained based on which a decision is made on whether to send an update request for this prefix. Available bandwidth is a valuable quantity that is measured and included in the computation of the performance score in some embodiments of the invention.

#### J. Differentiated Queues and Token Rates per Link

In some embodiments of the invention, token rates may differ on a per-link basis (which dictates the use of different queues for each link).

In some embodiments, the token rate may be tailored to total utilization. Lowly utilized links can afford relatively higher token rates without fear of  
25 overflow, whereas links close to saturation should be handled more carefully. Some embodiments of the invention provide one or more of the following modes of operation:

1. The default mode: the user specifies one token rate (and, optionally, a bucket size), shared equally among the prefixes updates destined to the different links.
2. The enhanced performance mode: the user specifies a minimum token rate (and, optionally, a bucket size). Depending on factors such as the total bandwidth utilization and the bandwidth of individual links, the prefix scheduler takes the initiative to function at a higher speed when possible, allowing better performance when it is not dangerous to do so.
3. The custom mode: in this case, the user can specify minimum and maximum token rates (and, optionally, bucket sizes), as well as conditions on when to move from a token rate to another. Using this custom mode, customers can tailor the prefix scheduler to their exact need.

#### K. Prefix Winner set *Re-computation*

Even though the priority queue is sized in such a way that the delay spent in the queue is minimized, there is still an order of magnitude between the time scale of the BGP world, at which level decisions are taken, and the physical world, in which edge stats and interface stats are measured. That is, even though the queuing delay is comparable to other delays involved in the process of changing a route, prefix performance across a given link or the utilization of a given link can change much more quickly. For example, a 2 second queuing delay could be appropriate in the BGP world, while 2 seconds can be enough for congestion to occur across a given link, or for the link utilization to go from 25% to 75%... For this reason, in some embodiments of the invention, the winner set is re-evaluated at the output of the priority queue.

#### L. Conclusion

The foregoing description of various embodiments of the invention has been presented for purposes of illustration and description. It is not intended to limit the invention to the precise forms disclosed. Many modifications and equivalent arrangements will be apparent.

## CLAIMS

What is claimed is:

1. A communications back-channel, for coordinating routing decisions, the communications back channel comprising:

5           a plurality of networking devices;

          a plurality of routing intelligence units, wherein each of the plurality of the plurality of routing intelligence units includes software for controlling a distinct subset of the plurality of networking devices, each of the plurality of routing intelligence units further including:

10                   one or more processes for controlling the distinct subset of networking devices; and

          one or more coordination processes for exchanging routing parameters with the plurality of routing intelligence units.

15           2. The communications back-channel of claim 1, wherein the one or more processes for controlling the distinct subset of networking devices are Border Gateway Protocol (BGP) sessions.

          3. The communications back-channel of claim 2, wherein each of the routing intelligence units is a route-reflector client.

20           4. The communications back-channel of claim 3, wherein each of the distinct subset of networking devices is a route reflector to the route reflector client.

          5. The communications back-channel of claim 1, wherein the one or more coordination process in each of the routing intelligence units includes BGP sessions.

25           6. The communications back-channel of claim 5, wherein the BGP sessions in the one or more coordination processes of each of the routing intelligence units includes:

at least one BGP process; and

at least one BGP stack, such that the at least one BGP stack exchanges routing parameters between the routing intelligence unit and the at least one BGP process, and the at least one BGP process exchanges routing parameters with the plurality of routing intelligence units.

7. The communications back-channel of claim 6, wherein the at least one BGP stack is a route reflector client, and the at least one BGP process is a route reflector.

8. The communications back-channel of claim 6, wherein the routing parameters include local path performance characteristics.

9. The communications back-channel of claim 6, wherein the routing parameters include performance scores for routes.

10. The communications back-channel of claim 9, wherein the performance scores are exchanged via a Local Preference field.

11. The communications back-channel of claim 1, further comprising:

a plurality of communication links directly coupling the plurality of routing intelligence units, wherein the plurality of communication links are dedicated exclusively for exchanging routing parameters between the plurality of routing intelligence units.

12. The communications back-channel of claim 11, wherein the plurality of communication links are at least partially comprised of physical links between the plurality of routing intelligence units.

13. The communications back-channel of claim 11, wherein the plurality of communication links are at least partially comprised of logical links between the plurality of routing intelligence units.

14. A method of exchanging routing parameters amongst a plurality of decision makers, each decision maker controlling a distinct subset of a plurality of routers, wherein the plurality of decision makers are in communication via a dedicated mesh, the method comprising:

5                   asserting a first plurality of preferred routes for a first plurality of prefixes to the subset of routers; and

                  concurrent with the asserting the first plurality of preferred routes, sending a plurality of local performance scores for the first plurality of routes to the plurality of decision makers via the dedicated mesh.

10           15. The method of claim 14, further comprising:

                  receiving a second plurality of routes for a second plurality of prefixes via the dedicated mesh.

                  16. The method of claim 15, further comprising:

15                   receiving a plurality of performance scores for the second plurality of routes.

                  17. The method of claim 16, wherein the plurality of performance scores are included in one or more Local Preferences fields in a BGP feed.

                  18. The method of claim 16, further comprising:

                  applying penalties to each of the plurality of performance scores.

20           19. The method of claim 14, wherein the asserting the first plurality of preferred routes is performed via a BGP feed to the subset of routers.

                  20. The method of claim 14, wherein the plurality of local performance scores are sent via a BGP feed to the dedicated mesh.

25           21. The method of claim 14, wherein the dedicated mesh is at least partially comprised of physical links between the plurality of decision makers.

22. The method of claim 14, wherein the dedicated mesh is at least partially comprised of logical links between the plurality of decision makers.



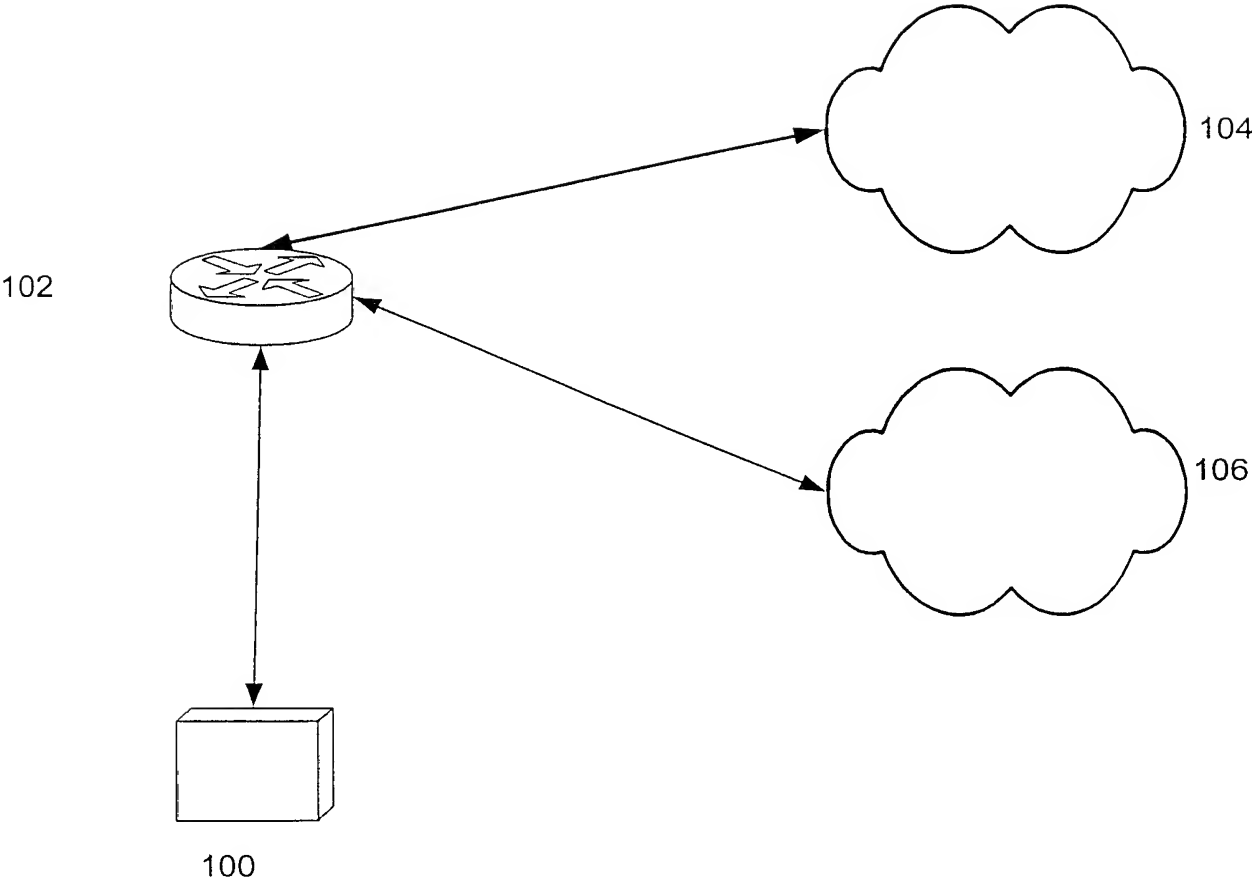


Figure 1

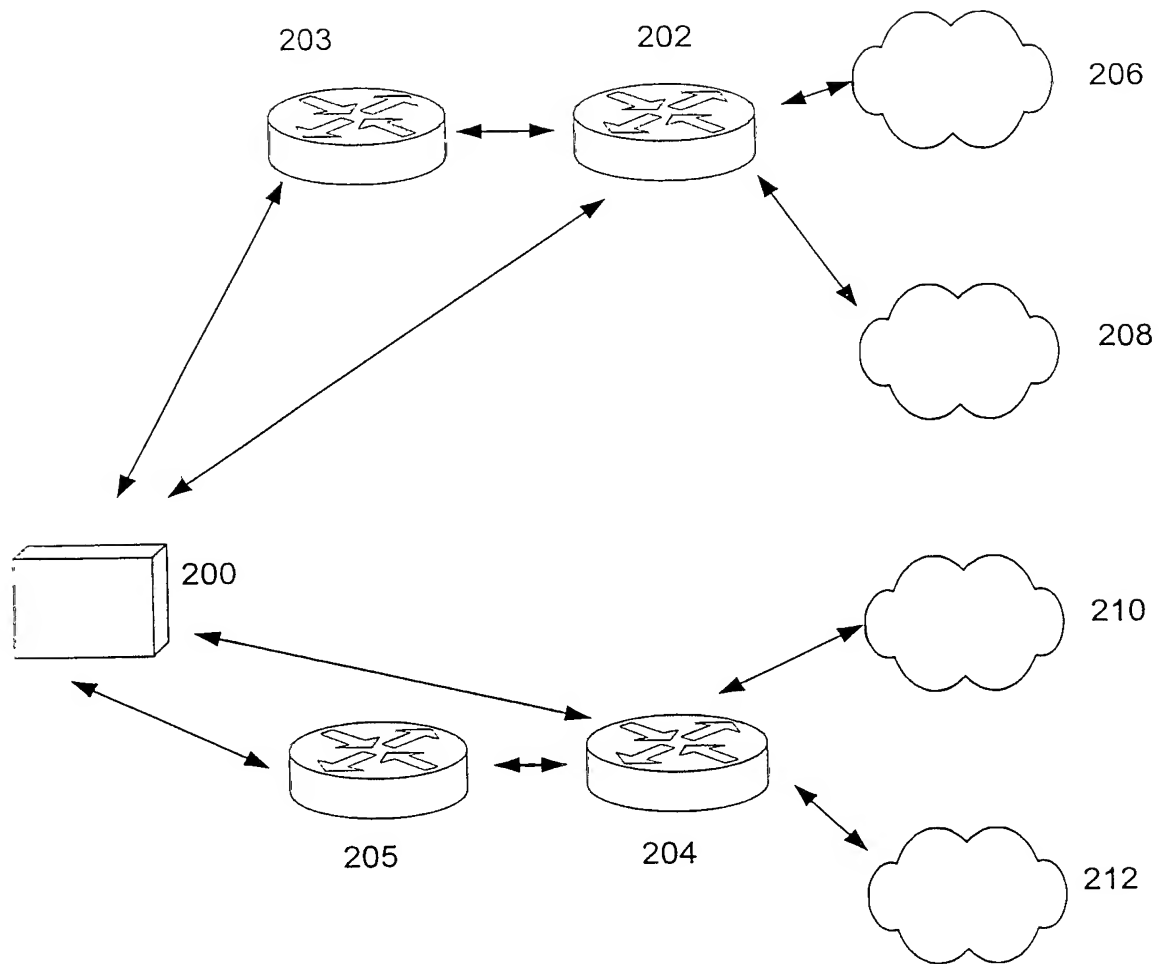
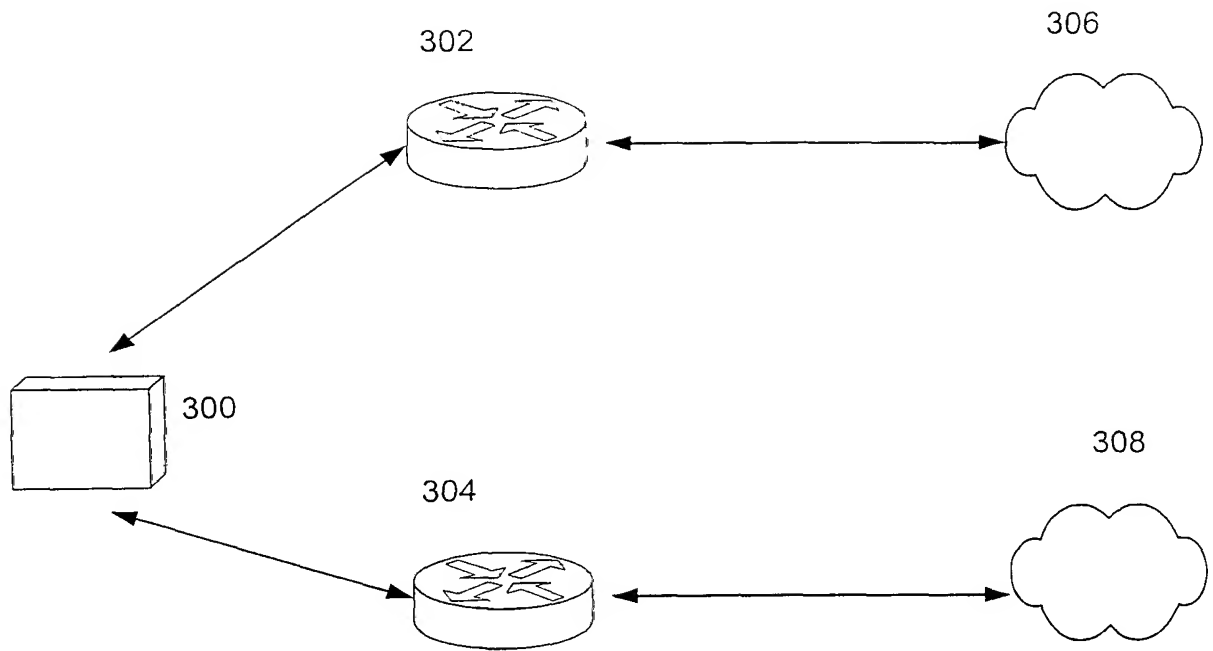


Figure 2

**Figure 3**

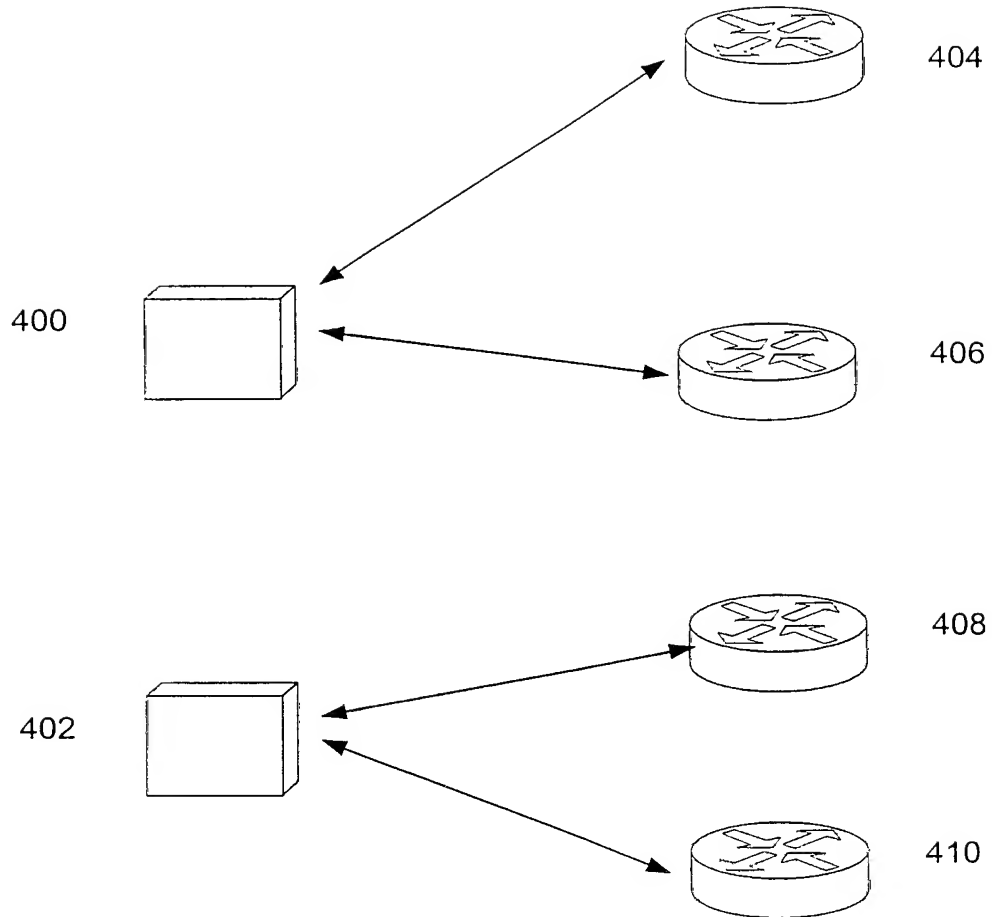
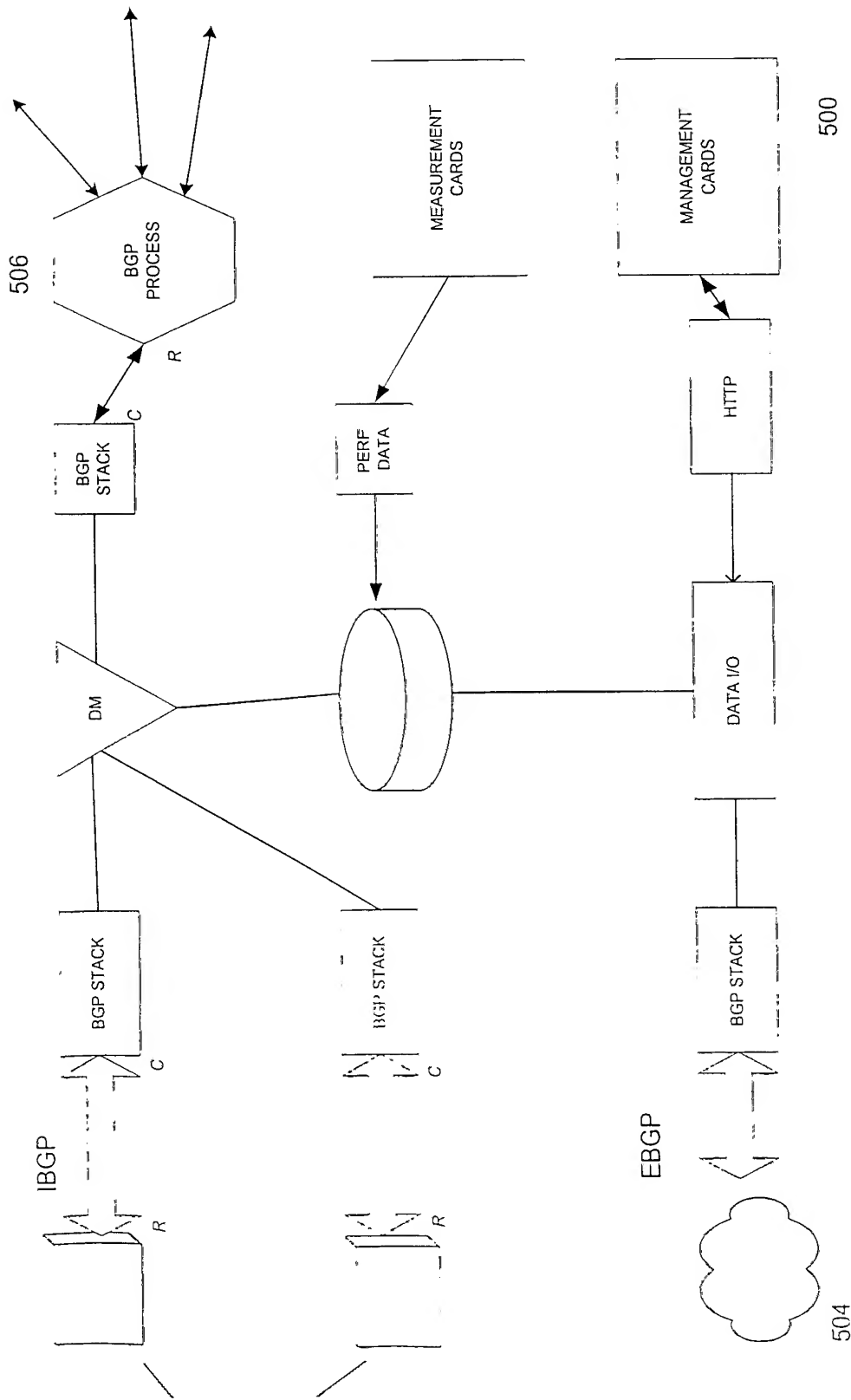


Figure 4



## Figure 5A

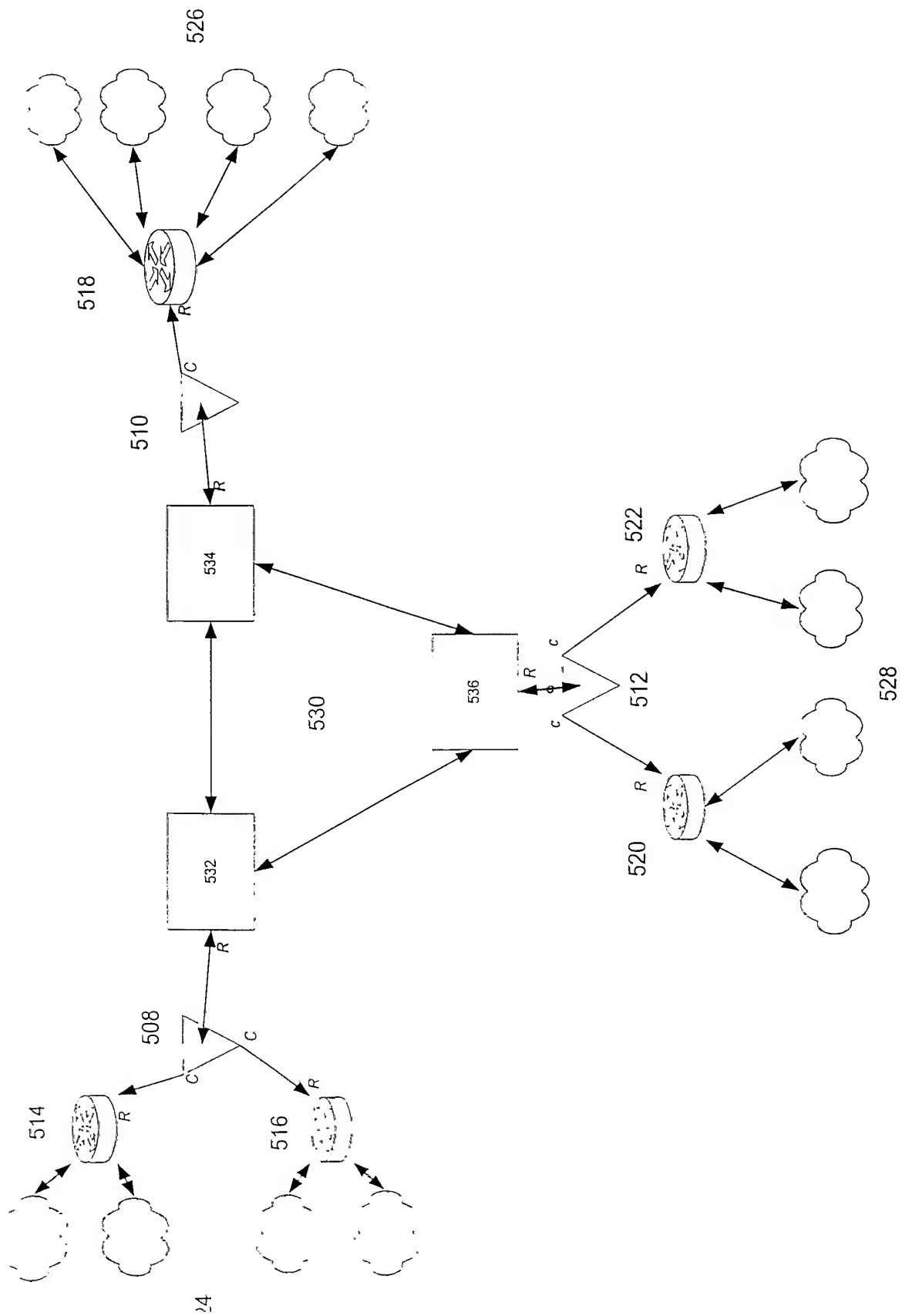


Figure 5B

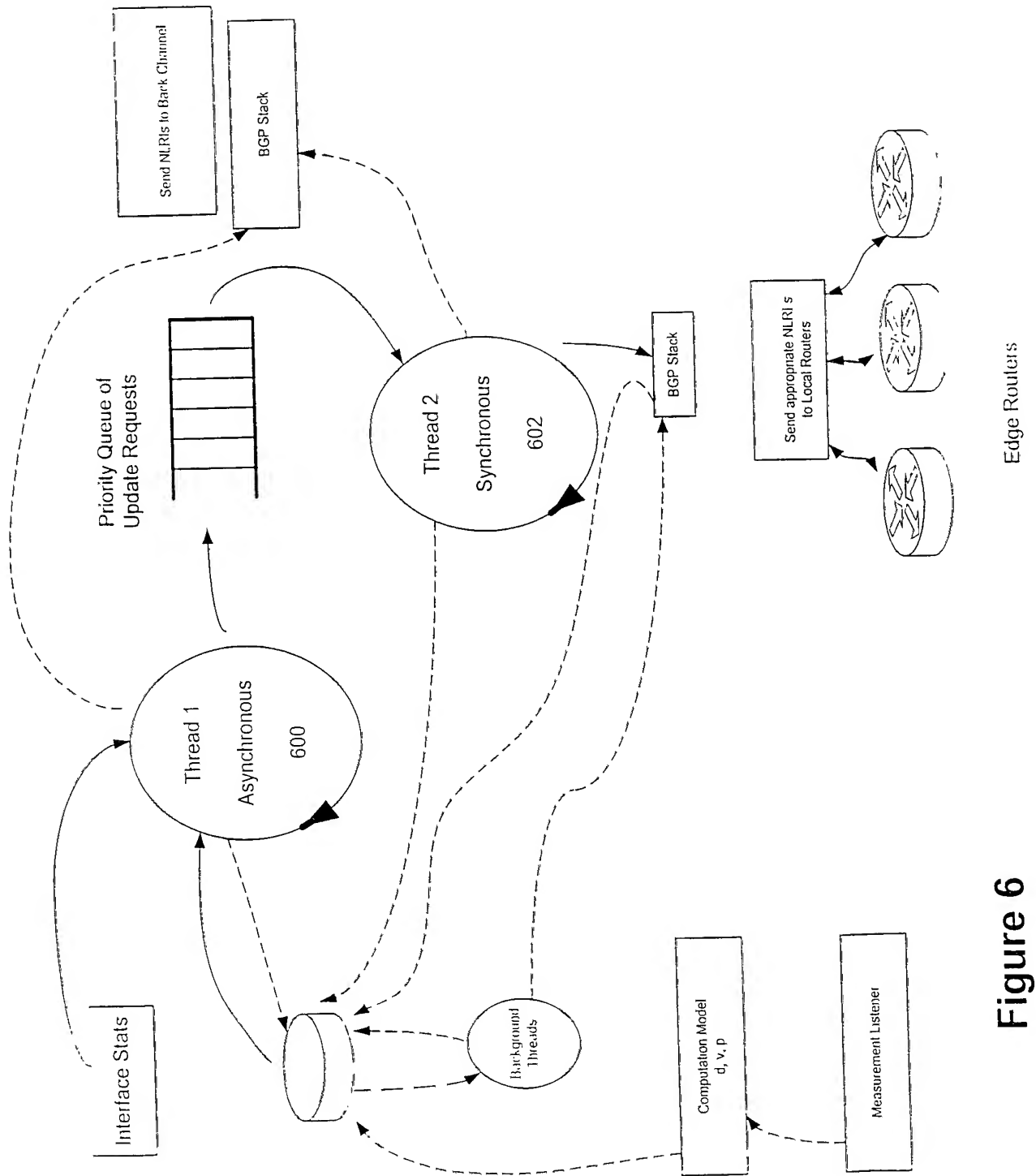


Figure 6

## INTERNATIONAL SEARCH REPORT

ational Application No

PCT/US 01/31259

A. CLASSIFICATION OF SUBJECT MATTER  
IPC 7 H04L12/56

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

IPC 7 H04L

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Electronic data base consulted during the international search (name of data base and, where practical, search terms used)

EPO-Internal, PAJ, INSPEC

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

| Category * | Citation of document, with indication, where appropriate, of the relevant passages   | Relevant to claim No. |
|------------|--|-----------------------|
| X          | J. YU: "Scalable Routing Design Principles"<br>RFC 2791 NETWORK WORKING GROUP,<br>31 July 2000 (2000-07-31), pages 1-24,<br>XP002191098<br>page 3, paragraph 4 - paragraph 11<br>page 6, paragraph 1 - paragraph 10<br>page 8, paragraph 7 -page 9, paragraph 6<br>page 9, paragraph 5 - paragraph 8<br>page 12, paragraph 9 -page 13, paragraph 4<br>page 15, paragraph 7 -page 16, paragraph 7 | 1,2,5,6               |
| A          | -----  | 3,7-9,<br>11,14       |

☐ Further documents are listed in the continuation of box C.

☐ Patent family members are listed in annex.

## \* Special categories of cited documents :

- \*A\* document defining the general state of the art which is not considered to be of particular relevance
- \*E\* earlier document but published on or after the international filing date
- \*L\* document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)
- \*O\* document referring to an oral disclosure, use, exhibition or other means
- \*P\* document published prior to the international filing date but later than the priority date claimed

- \*T\* later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention
- \*X\* document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone
- \*Y\* document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art.
- \*Z\* document member of the same patent family

Date of the actual completion of the international search

4 March 2002

Date of mailing of the international search report

22/03/2002

Name and mailing address of the ISA

European Patent Office, P.B. 5818 Patentlaan 2  
NL - 2280 HV Rijswijk  
Tel (+31-70) 340-2040, Tx 31 651 epo nl,  
Fax: (+31-70) 340-3016

Authorized officer

Brichau, G